

# **SANAL ARAŞTIRMA ORTAMLARI ve AÇIK VERİLER**

Bülent Karasözen, ODTÜ

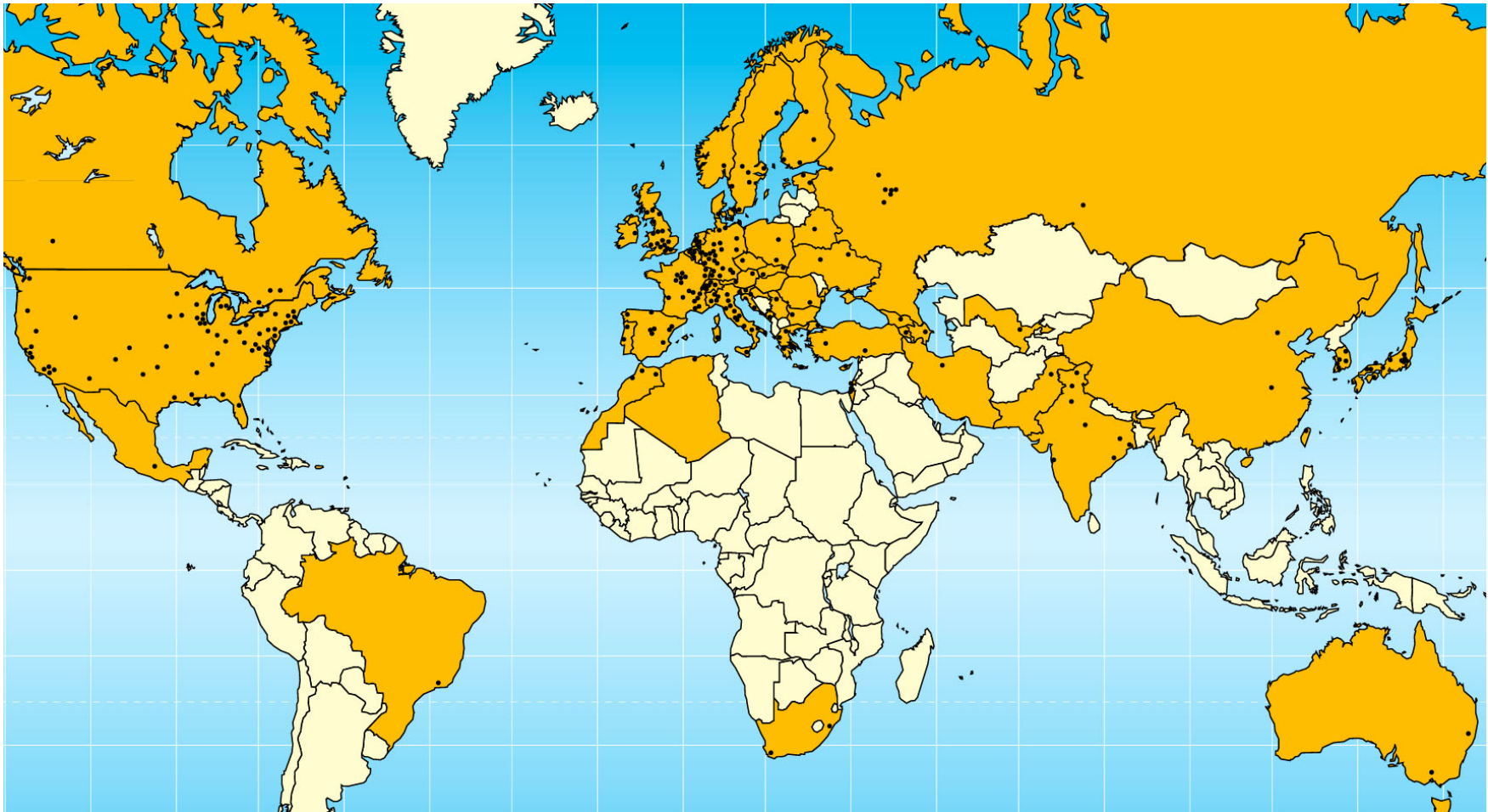
INER-TR 05

9-11 Aralık 2005, Bahçeşehir Üniversitesi

# E-bilim

- Bilim ve teknolojide yeni problemlerin çözümü için giderek artan küresel işbirliği ve kaynakların ortak kullanımı
- 1990, Tim Barnes Lee , CERN, World Wide Web
- 1990, J.C.R. Licklider, ARPANET

# CERN kullanıcıları



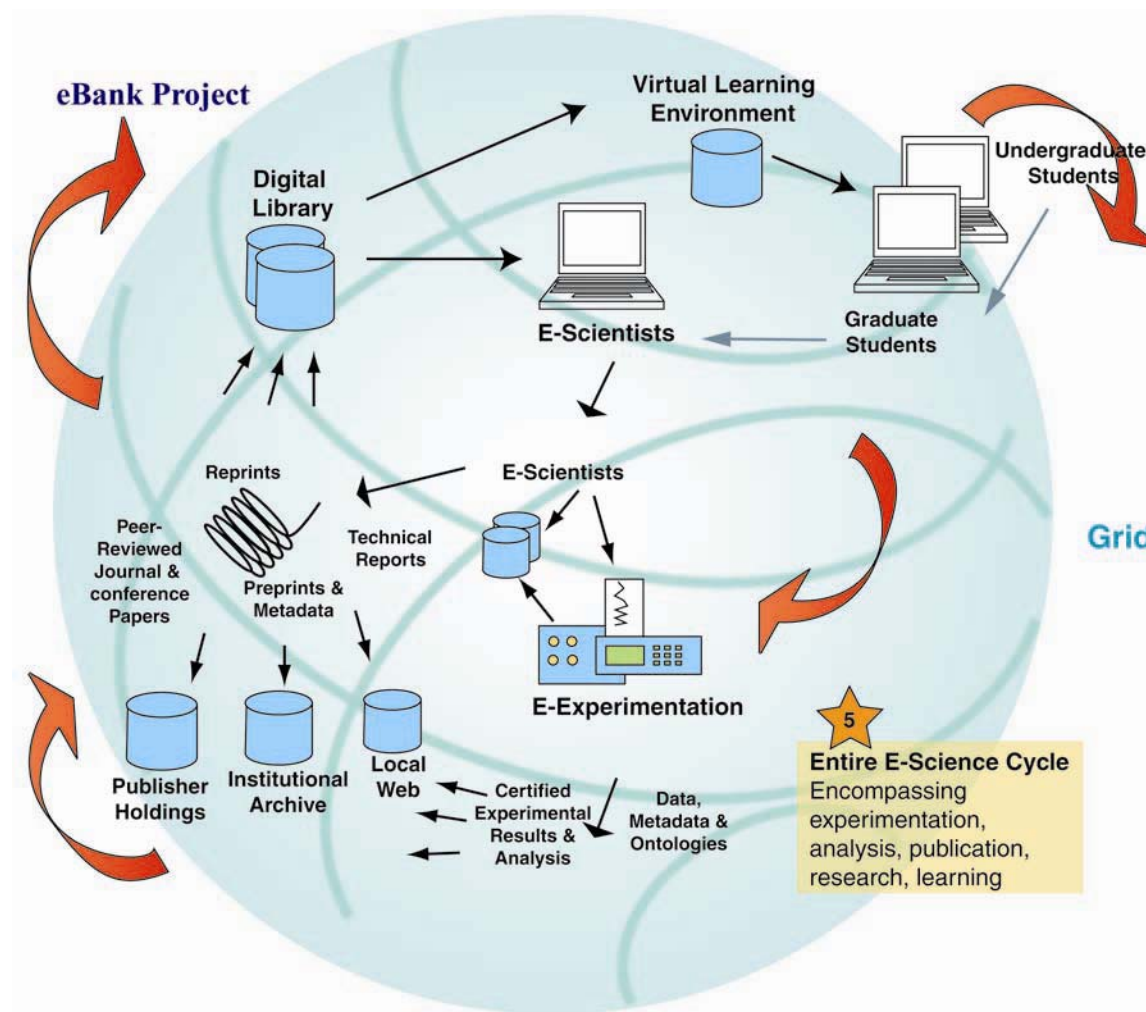
**Avrupa: 267 enstitü, 4603 kullanıcı**  
**Diğer ülkeler: 208 enstitü, 1632 kullanıcı**

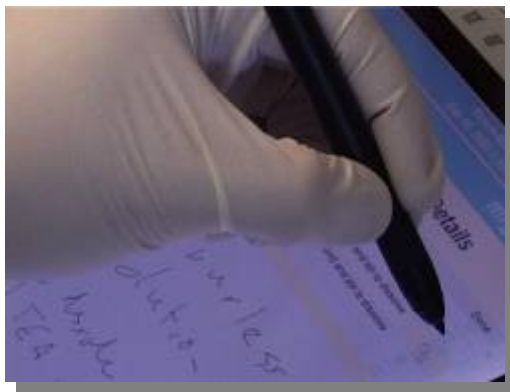
# **E-bilim projelerinden bazı örnekler**

- **Parçaçık fiziği**
  - Deneylerde elde edilen verilerin küresel düzeyde paylaşımı ve simülasyonlarda kullanımı
- **Astronomi**
  - Teleskoplardan elde edilen verilere dayalı ‘sanal gözlem evleri’ oluşturulması
- **Kimya**
  - Deney aletlerinin uzaktan kontrolü ve elektronik labrotuvar kitapçığı
- **Bioinformatik**
  - Verilerin bütünleştirilmesi, veri akış planları, veriye dayalı bilgi öğretimi
- **Sağlık**
  - normalize edilmiş mammogramların paylaşımı
- **Çevre**
  - İklim modellemesi

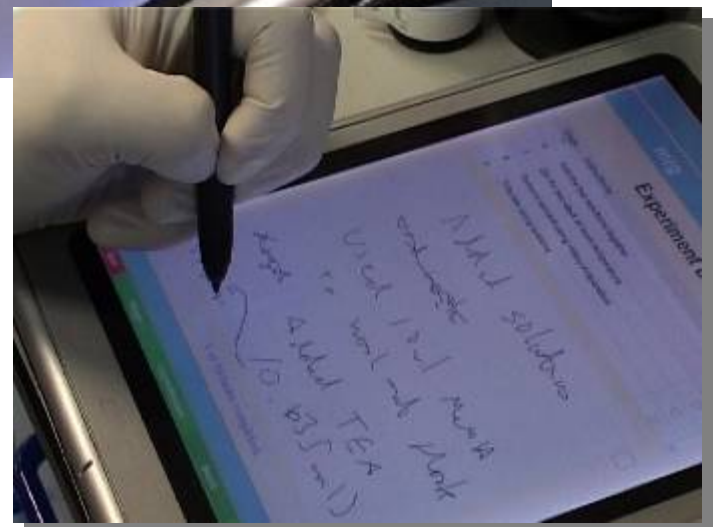
# E-Bilimi itekliyen ana etkenler

- Büyük çaptaki sayıları az pahalı labratuvar, super bilgisayar ve veri depolarına erişim:
  - CERN LHC
- Açık kaynaklı, kaliteli 'grid middlewaere'nin oluşturulması:
- - OMII, NMI, C-Omega
- Araştırmacıların karşısına çıkan veri seli:
  - Parçacık fiziği, astronomi, bioinformatik
- Açık erişim hareketi:
  - Bilimsel yayınlara ve verilere açık erişim





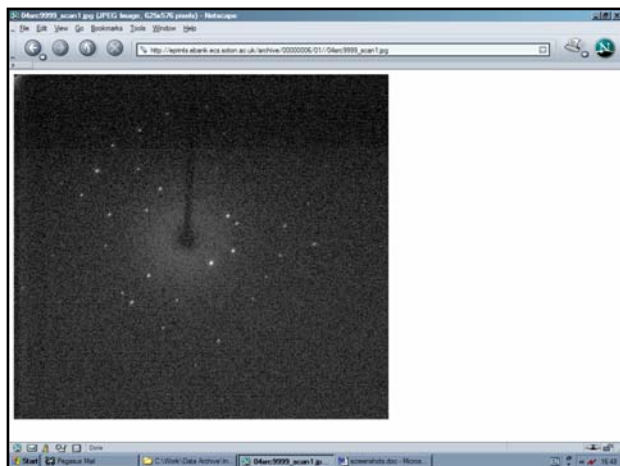
digital lab book





# Crystallographic e-Prints

➤ Direct Access to Raw Data from scientific papers



EBank Southampton - N-(5-Chloro-pyridin-2-yl)-4-fluoro-benzamide C12H8ClF N2 O

Deposited By: Christopher Gutteridge  
Deposited On: 26 February 2004

data

- 04src9999\_nonius-config.py (1179)
- 04src9999\_scan1.jpg (127447)
- 04src9999\_scan2.jpg (126927)
- 04src9999\_drx (849)
- 04src9999\_non (2645)
- 04src9999\_mmt (287)
- 04src9999\_scheme.jpg (3069)
- extracted\_cml (3325)
- extracted\_data.txt (489)

Chemical Formula: C12 H8 Cl F N2 O  
CFOM: 0.0366  
Cell Angle Alpha: 77.641(4)  
Symmetry Cell Setting: triclinic  
Symmetry Space Group Name H-M: P-1  
Cell Angle Gamma: 86.374(6)  
Cell Angle Beta: 80.643(6)  
Refine LS R Factor All: 0.1079  
Refine LS WR Factor GT: 0.1091  
Refine LS WR Factor Ref: 0.1292  
Cell Length A: 5.2061(3)  
Cell Length B: 10.2615(11)  
Diffraction Ambient Temperature: 120(2)  
Refine LS R Factor GT: 0.0531  
Cell Length C: 10.6118(10)  
Exptl Crystal Description: plate

Archive Staff Only: [edit this record](#)

Contact Information

EPSRC National Crystallography Service  
Data Collection Summary kced1 (dellboy)

Summary report for Directory: diska/03hms003  
Report generated Mar 19, 2003, 17:01:34

Unit cell

6358 reflections with 2.91° delta (2° 43" resolution between 7.00 Å and 0.77 Å) were used for unit cell refinement

Symmetry used in scaling: P1

a (Angstrom)	5.2064 ± 0.0003
b (Angstrom)	10.2621 ± 0.0011
c (Angstrom)	10.6123 ± 0.0010
alpha (°)	77.643 ± 0.004
beta (°)	80.636 ± 0.006
gamma (°)	86.374 ± 0.006
Volume (Å³)	546.25 ± 0.08
Mosadity (°)	0.673 ± 0.004

Raw data sets can be very large and these are stored at National Datastore using SRB server



# NSF 'Atkins' siber altyapı raporu:

- Bir çok bilim dalında en son buluşlara ve yeniliklere erişim WEB üzerinden gerçekleşmekte
- Yüzlerce ve binlerce terabyte bilimsel verinin arşivlenmesi ve erişime sunulması, bilim ve teknolojinin ilerlemesi için artık vazgeçilemez bir gereksinim

## MIT'nin DSpace vizyonu:

- Araştırmacılar tarafından üretilen makaleler, raporlar, veri setleri, deney sonuçları genellikle kişisel veya bölüm WEB sayfalarında saklanmakta. Araştırmacıların ayrılması veya bölümlerin zaman için de değişime uğramasıyla kaybolabilmekte

# Veri Seli

- Büyük bir roman: 1 Mbyte
- İncil: 5 Mbytes
- Bir Mozart semfonisi (sıkıştırılmış): 10 Mbytes
- OED on CD: 500 Mbytes
- Dijital film (sıkıştırılmış): 10 Gbytes
- Hakemli bir derginin bir yıllığı (~20k dergi; ~2M makale): 1 Tbyte
- Library of Congress: 20 Tbytes
- İnternet arşivi (10 milyon sayfa) ( 1996 – 2002 arası): 100 Tbytes
- Yıllık basılı, film, optik ve manyetik medya üretimi: 1500 Pbytes

**Basılı kaynaklar, toplam saklanan bilginin sadece 0.003% 'ünü oluşturmakta**

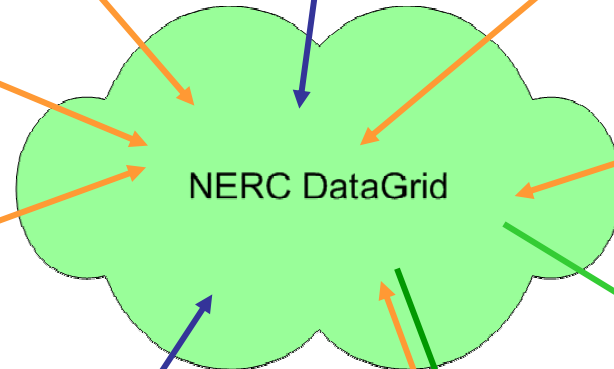
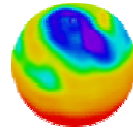
**Önümüzdeki beş yılda e-bilim projelerinden, insanlık tarihindeki tüm verilerin toplamından kat, kat fazla petabyte büyüklüklerinde veri üretilecek**

# Grid, middleware

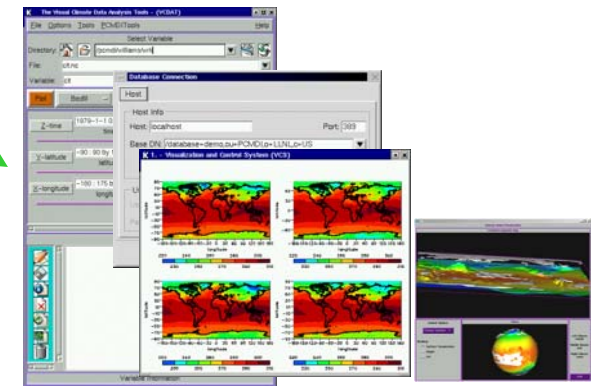
- Grid: birbirinden farklı çeşitli bilgisayarların ortak kullanımının sağlandığı ortam
- Büyük çapta çok disiplinli, araştırmacıların farklı yerlerde olduğu ortak projelerin statik WEB sayfaları aracılığıyla gerçekleştirilmesi mümkün değil
- Middleware: bilgisayar ağıyla, uygulamalar arasında iletişimi sağlayan yazılım

# Complexity + Volume + Remote Access = Grid Challenge

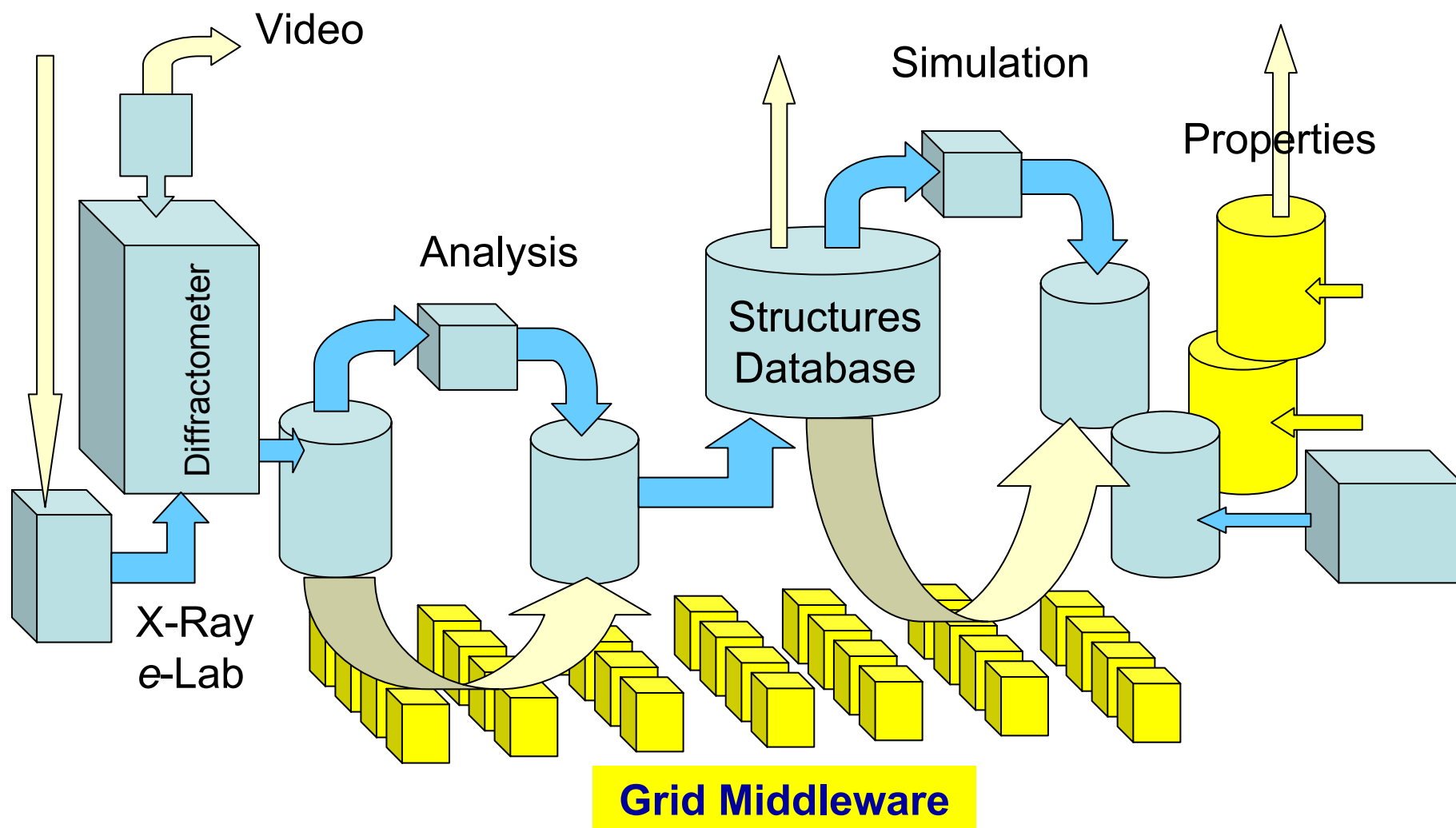
British Atmospheric  
Data Centre



British  
Oceanographic Data  
Centre



# Comb-e-Chem Project



# **Açık veriler için metadata standartları**

San Diego metadata şeması:

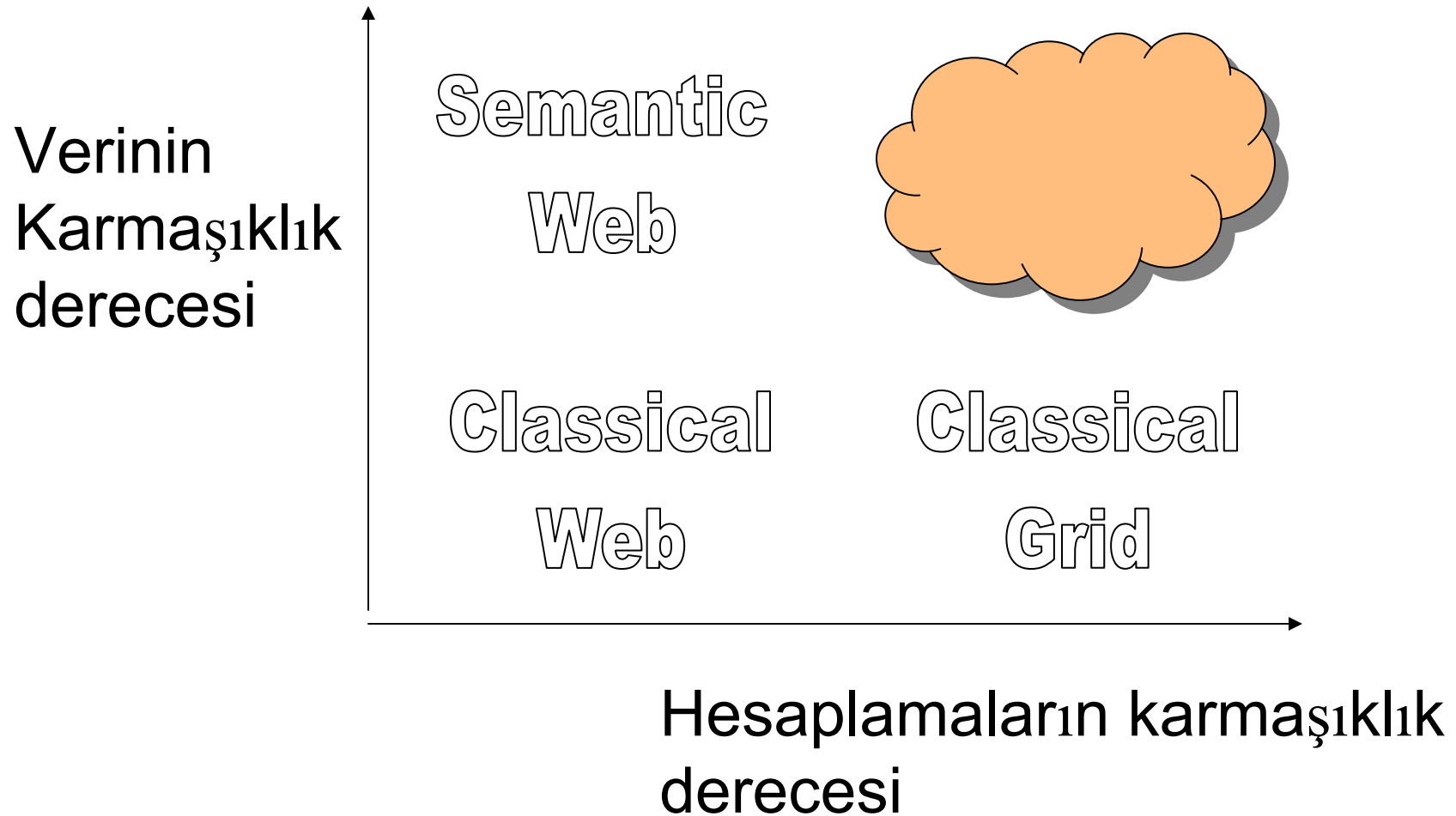
Verilerin saklaması ve erişimi için gerekli sistem metadata

Verinin özelliklerini içeren Dublin Core'a dayalı metadata

Kullanıcının erişim bilgilerini ve haklarını içeren metadata

Araştırma alanına özgü bilgileri içeren metadata

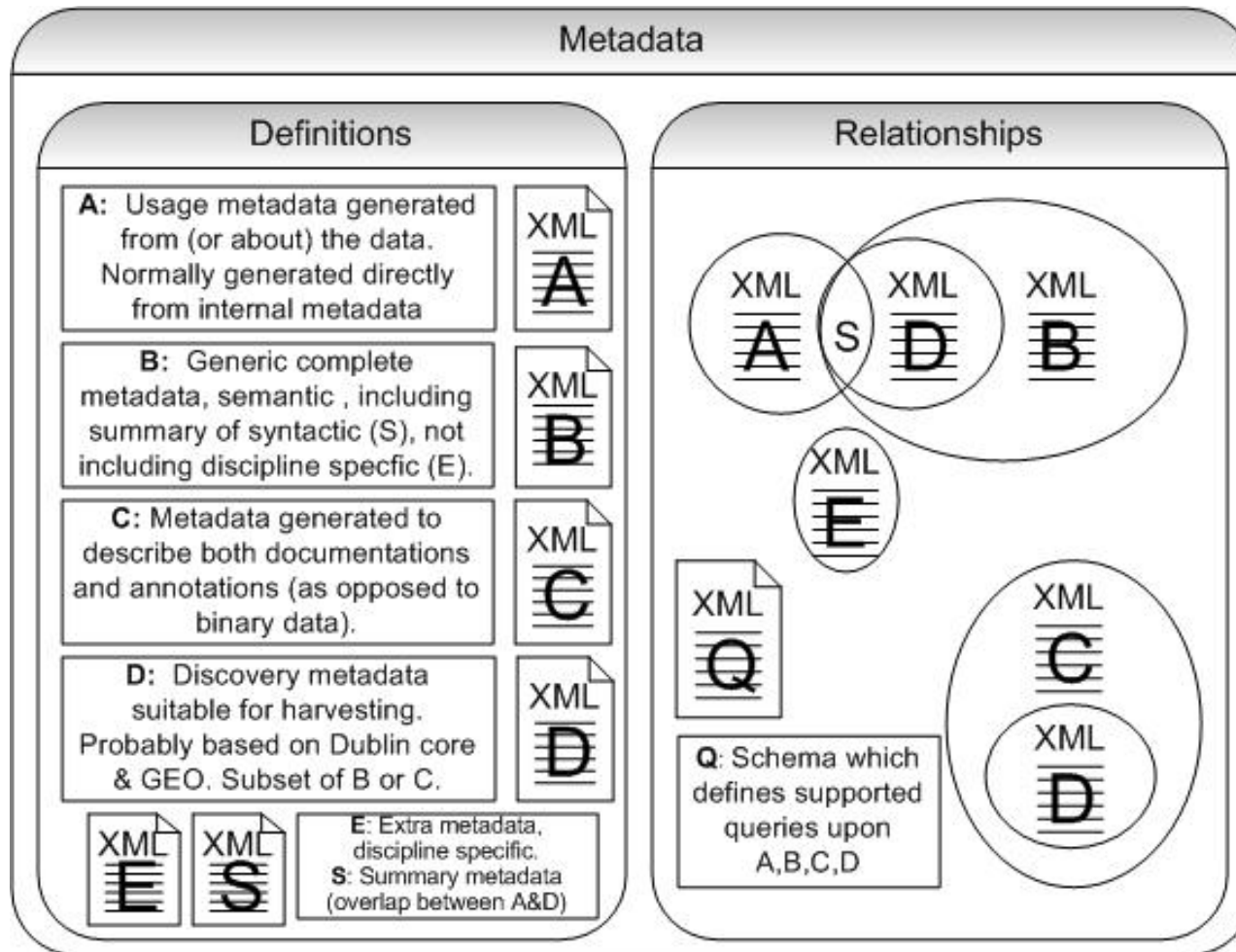
# veriden bilgiye semantik 'grid'





# NERC Data Grid

## Metadata Taxonomy



# **E-bilim ve kütüphaneler**

Önümüzdeki on yıl içerisinde üniversite kütüphanelerinin en önemli görevi, üniversitede üretilen bilimsel yayınların saklanması ve açık erişime sunulması olacak

Bilimsel yayınların yanı sıra, araştırmalarla ilgili verilerin de erişime sunulması ve saklanması önem kazanmakta

# Verilerin bakımı

- Hastahanelerdeki milyonlarca filmin arşivlenmesi, işlenmesi
- Farklı dillerdeki yazılımların ilerde de kullanılabilirliği
- Zaman içinde değişen verilere uygun metadata standartları geliştirilmesi
- Farklı formatlardaki veriler arasında iletişim sağlanması (interoperability)